# UK Federation Plans and Initiatives
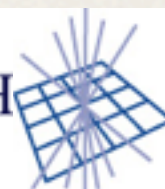
Wahid Bhimji

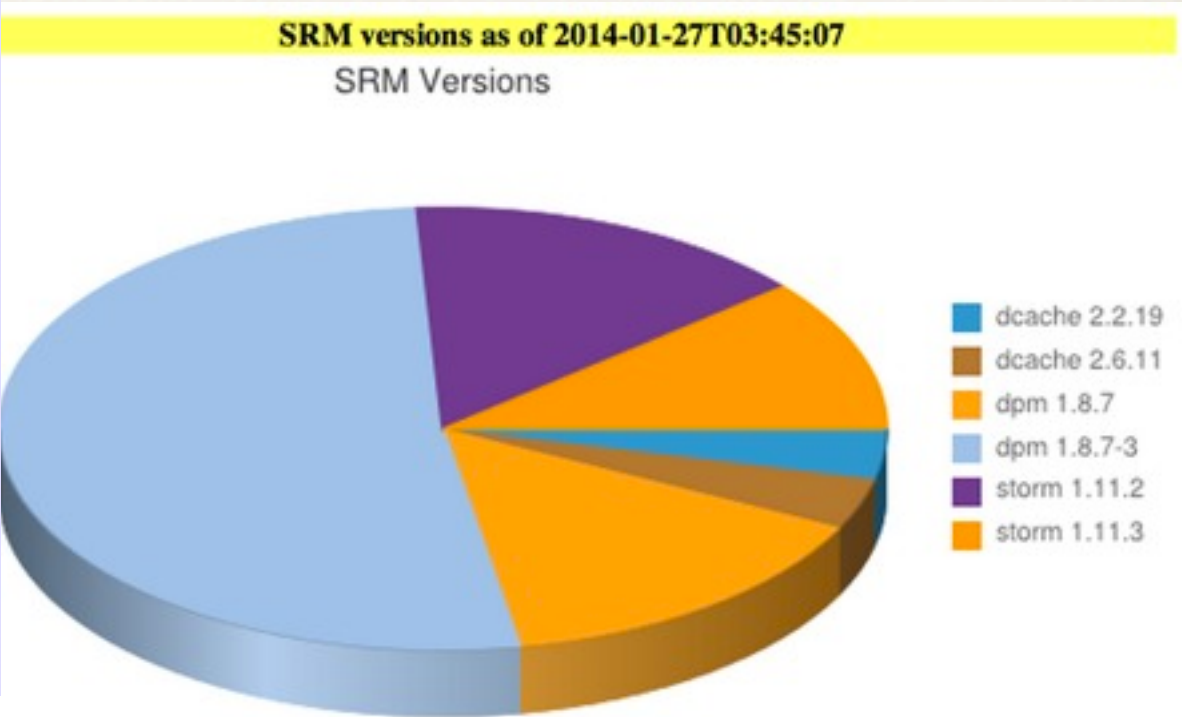THE UNIVERSITY *of* EDINBURGH

GridPP
UK Computing for Particle Physics

# Data federations: a UK perspective

* ATLAS and CMS are using Federated Storage in production now and more so in Run 2:

    * The UK is heavily involved and largely deployed

* We need to consider impact on our infrastructure. So need to :

    * Compare experiences across experiments e.g. job profile and needs.

    * Get information on (projected) use: TEG said federation traffic would be <~10% of bandwidth; now see >10%, total expected traffic would also be interesting.

    * Conduct our own infrastructure tests and evaluate monitoring

* HTTP / DAV - promising ideas that will be realized during Run 2. Sites currently having to run both, need to know how this heterogeneous (xrd/http) landscape will evolve ...

# The UK (1707-2014)



* Tier 1 (RAL) runs CASTOR (currently)

* Many Tier 2 sites, some small, grouped in "federated" T2s

* Most Tier2 sites run DPM

* Most sites support more than one experiment though focus on CMS (IC,RALPP, Brunel, Bris) or ATLAS (the rest). Bham is also an ALICE site.
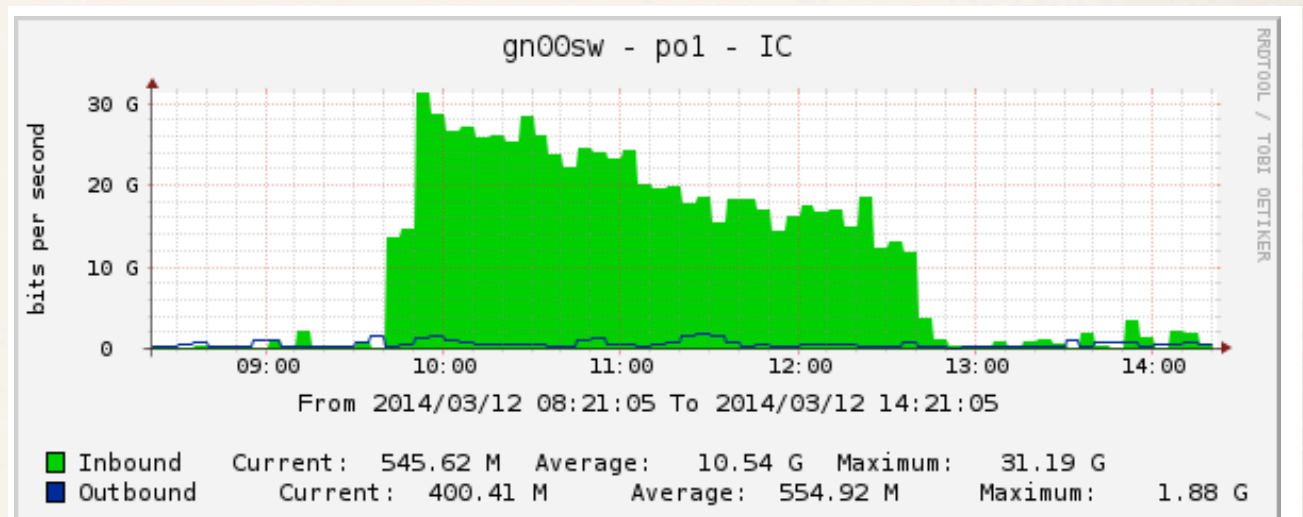


SRM versions as of 2014-01-27T03:45:07

SRM Versions

- dcache 2.2.19
- dcache 2.6.11
- dpm 1.8.7
- dpm 1.8.7-3
- storm 1.11.2
- storm 1.11.3

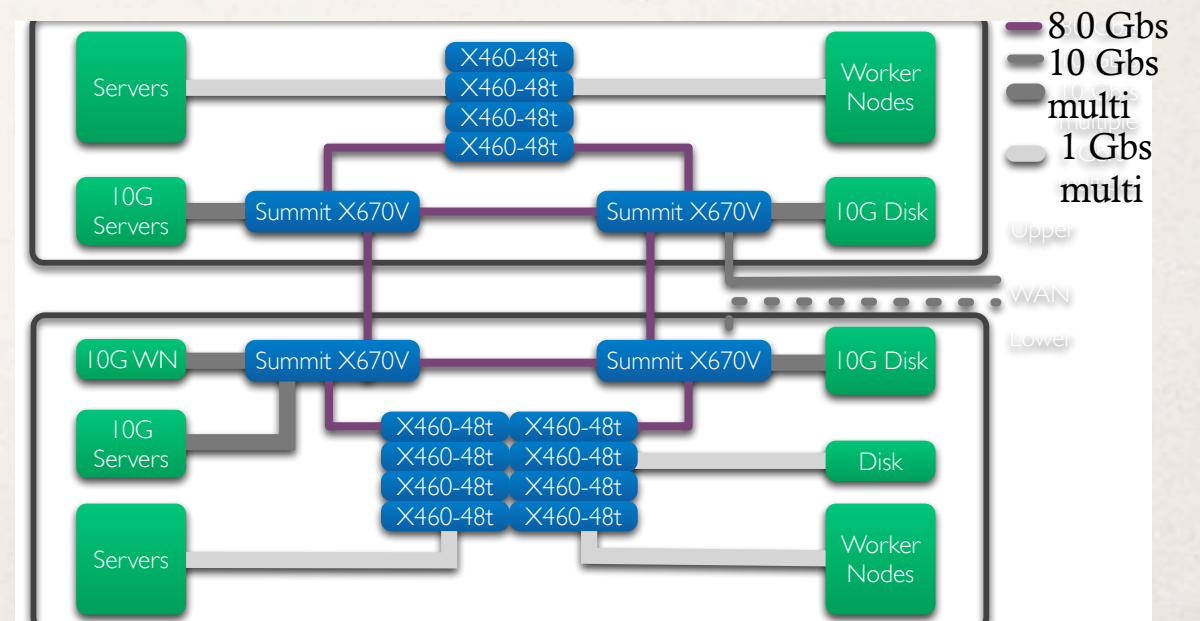| Site | Type | Version | Disk [TB] |
|------|------|---------|-----------|
| UKI-LT2-Brunel | dpm | 1.8.7-3 | 500 |
| UKI-LT2-IC-HEP | dcache | 2.6.11 | 2000 |
| UKI-LT2-QMUL | storm | 1.11.3 | 1700 |
| UKI-LT2-RHUL | dpm | 1.8.7-3 | 600 |
| UKI-LT2-UCL-HEP | dpm | 1.8.7-3 | 190 |
| UKI-NORTHGRID-LANCS-HEP | dpm | 1.8.7-3 | 1000 |
| UKI-NORTHGRID-LIV-HEP | dpm | 1.8.7-3 | 550 |
| UKI-NORTHGRID-MAN-HEP | dpm | 1.8.7-3 | 1000 |
| UKI-NORTHGRID-SHEF-HEP | dpm | 1.8.7-3 | 400 |
| UKI-SCOTGRID-DURHAM | dpm | 1.8.7-3 | 50 |
| UKI-SCOTGRID-ECDF | dpm | 1.8.7-3 | 350 |
| UKI-SCOTGRID-GLASGOW | dpm | 1.8.7-3 | 1300 |
| UKI-SOUTHGRID-BHAM-HEP | dpm | 1.8.7-3 | 300 |
| UKI-SOUTHGRID-BRIS-HEP | storm | 1.11.2 | 100 |
| UKI-SOUTHGRID-CAM-HEP | dpm | 1.8.7-3 | 300 |
| UKI-SOUTHGRID-OX-HEP | dpm | 1.8.7-3 | 650 |
| UKI-SOUTHGRID-RALPP | dcache | 2.2.19 | 1250 |
| UKI-SOUTHGRID-SUSX | storm | 1.11.2 | 50 |

# Networking

* Upgrades at Tier 1 (~now): will have resilient (dual) 40 Gbit/s connection to Janet-6 (the Tier2s)

* Most Tier2 sites have (shared or dedicated) 10 Gbit links to WAN and internally to disk servers. Planning for some 20 Gbit

* But some have less. Also external WN traffic may go through NAT

* Not at US levels, but think per "federated" T2 or per TB ...

Our best connected Tier2 (IC):



Another lucky one (Glasgow):

# ATLAS UK - status and testing

✤ Tier1 and almost all larger ATLAS T2s are in FAX (one remaining site also supports ALICE ...). Decent availability:

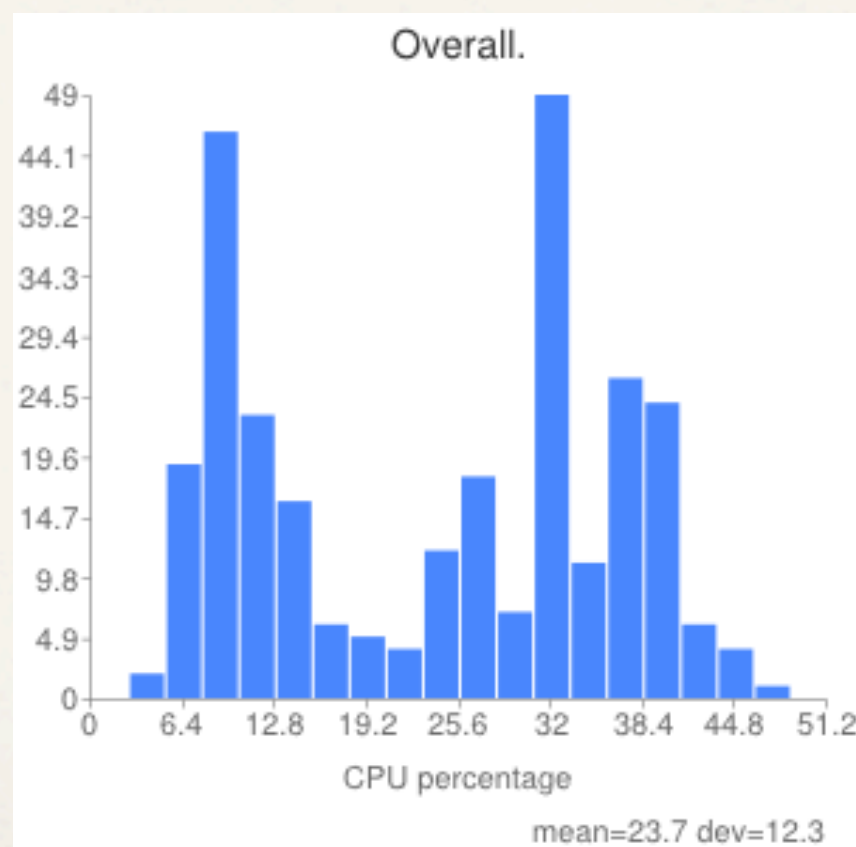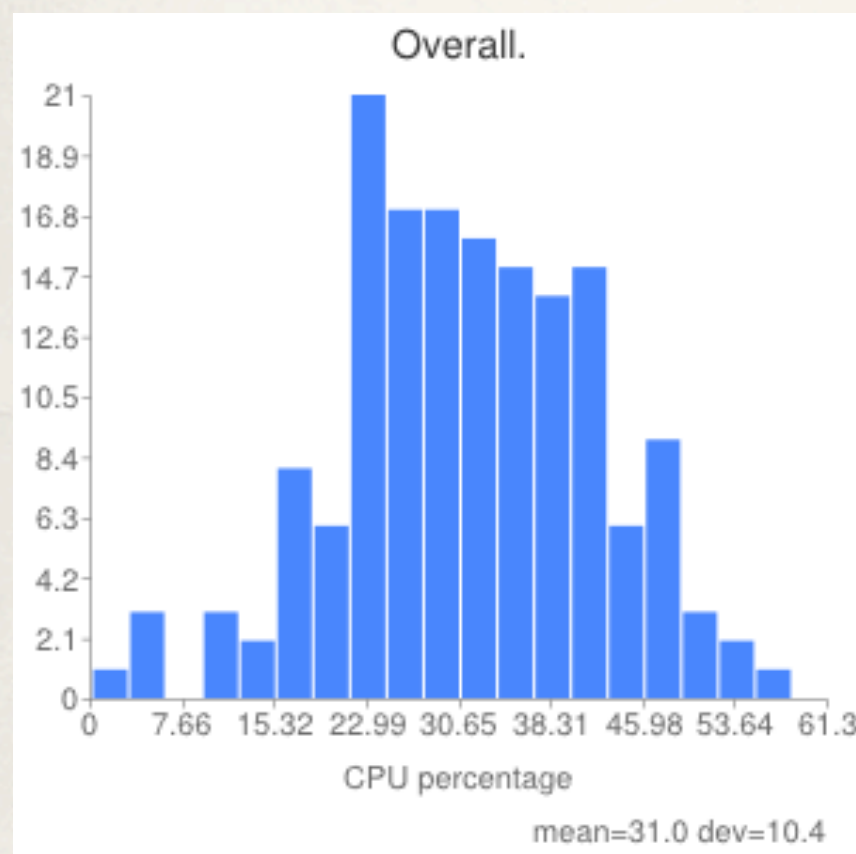| | | | |
|---|---|---|---|
| UKI-LT2-QMUL | OK | OK | OK |
| UKI-LT2-RHUL | noDirect | NoUpstreamRedirection | NoFirstLevelRedirection |
| UKI-NORTHGRID-LANCS-HEP | OK | OK | OK |
| UKI-NORTHGRID-LIV-HEP | OK | OK | OK |
| UKI-NORTHGRID-MAN-HEP | OK | OK | OK |
| UKI-NORTHGRID-SHEF-HEP | OK | OK | OK |
| UKI-SCOTGRID-ECDF | OK | OK | OK |
| UKI-SCOTGRID-GLASGOW | OK | OK | OK |
| UKI-SOUTHGRID-CAM-HEP | OK | OK | OK |
| UKI-SOUTHGRID-OX-HEP | OK | OK | OK |

✤ Started to perform "stress" tests (similar to those done in US) on "standard" (user.flegger.*) FAX (SMWZ) D3PD dataset:

  ✤ Direct reading in ROOT script (10% of events, 30 MB TTreeCache)
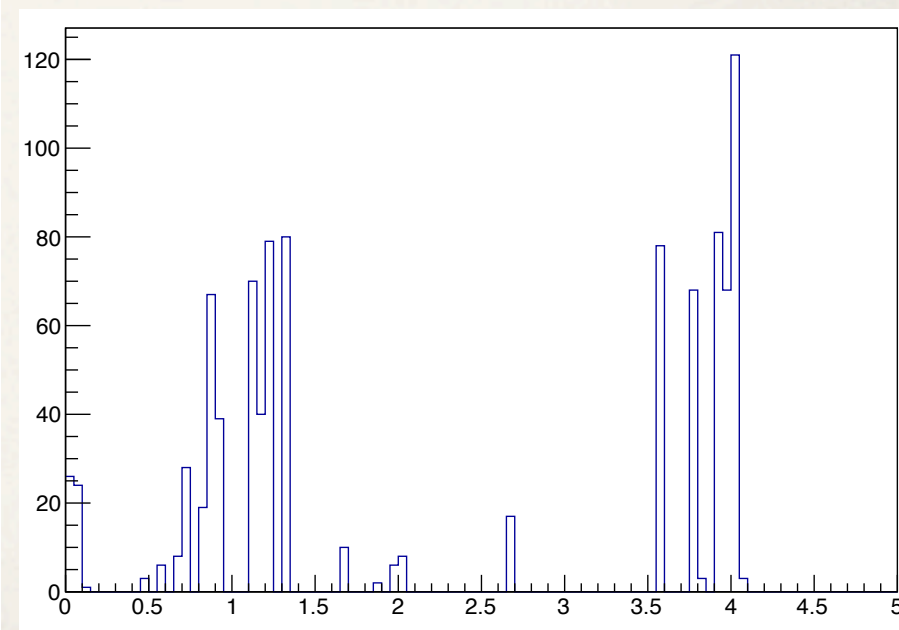
  ✤ Real H->WW analysis via Hammercloud

# ATLAS UK testing

Initial results from DPM sites (100-200 H->WW jobs): performance is fine
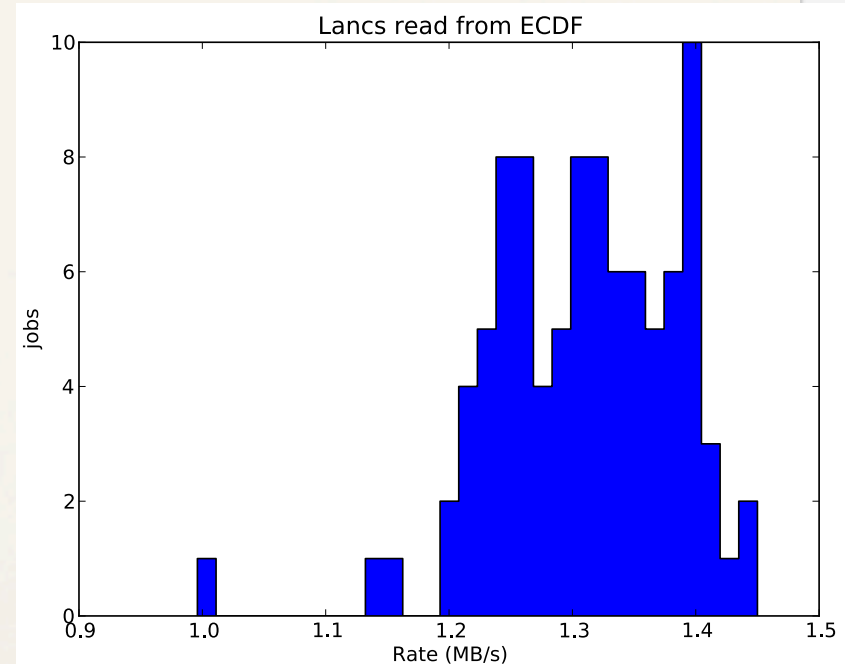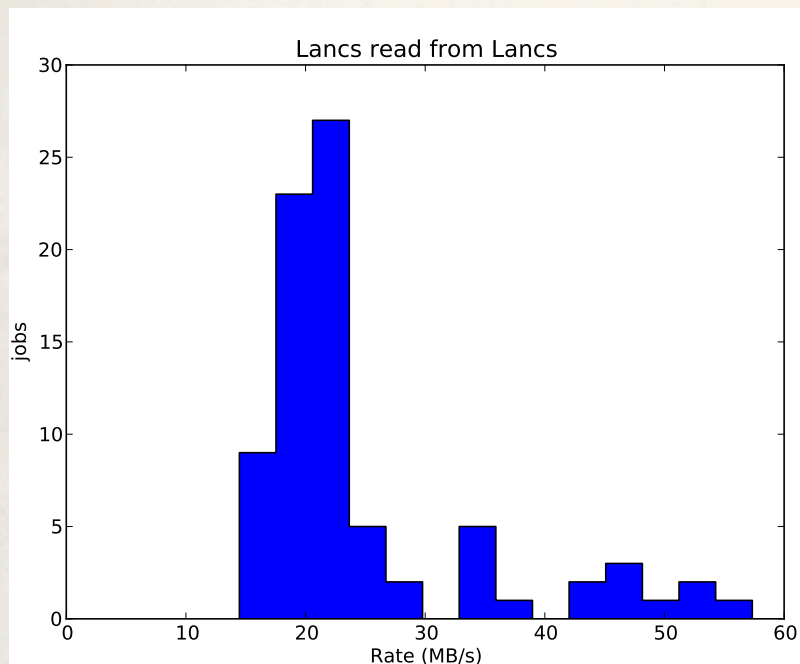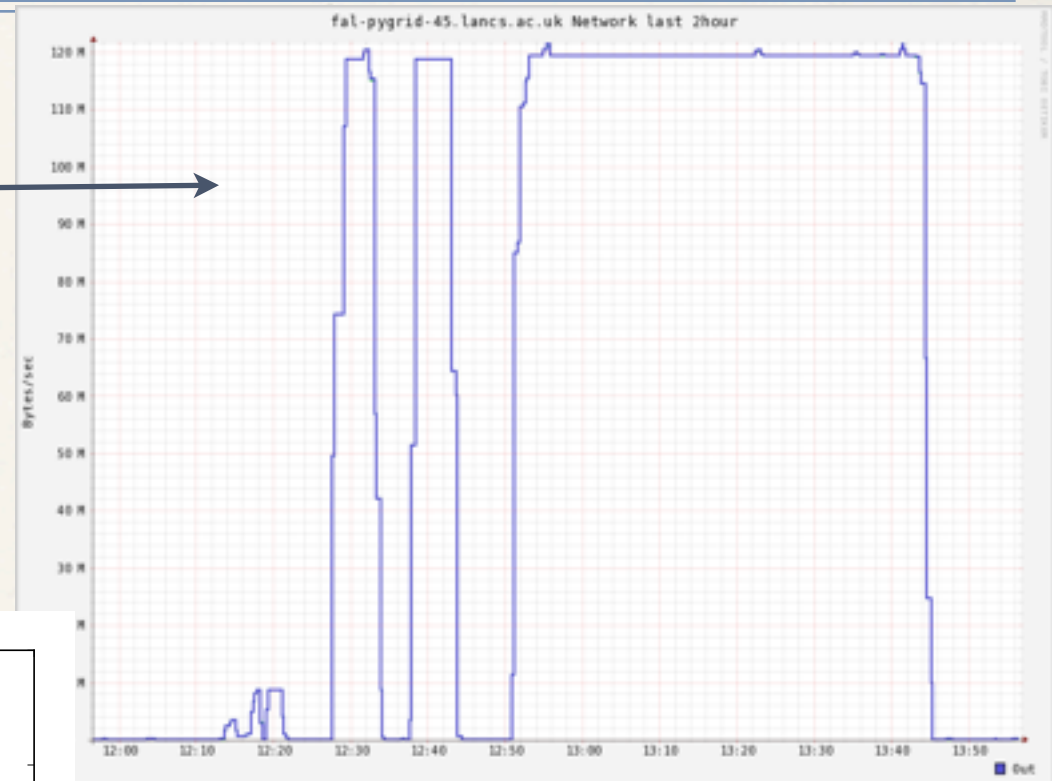
Oxford local read:    Oxford read from ECDF:

Data rate from xrootd monitoring (MB/s):

# Finding bottlenecks...

✤ 1Gig NAT easily saturated
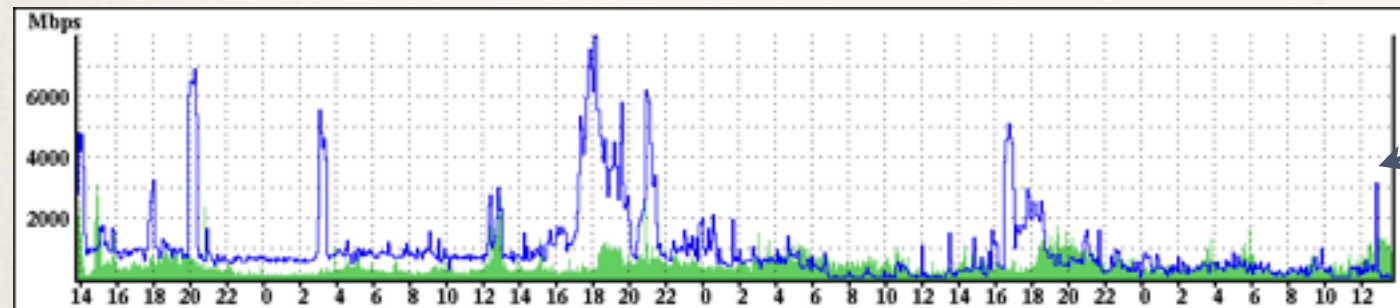(running 100 jobs with simple
ROOT script on local batch)


fal-pygrid-45.lancs.ac.uk Network last 2hour


Lancs read from Lancs


Lancs read from ECDF

Local > 20 MB/s read,
decreases externally to
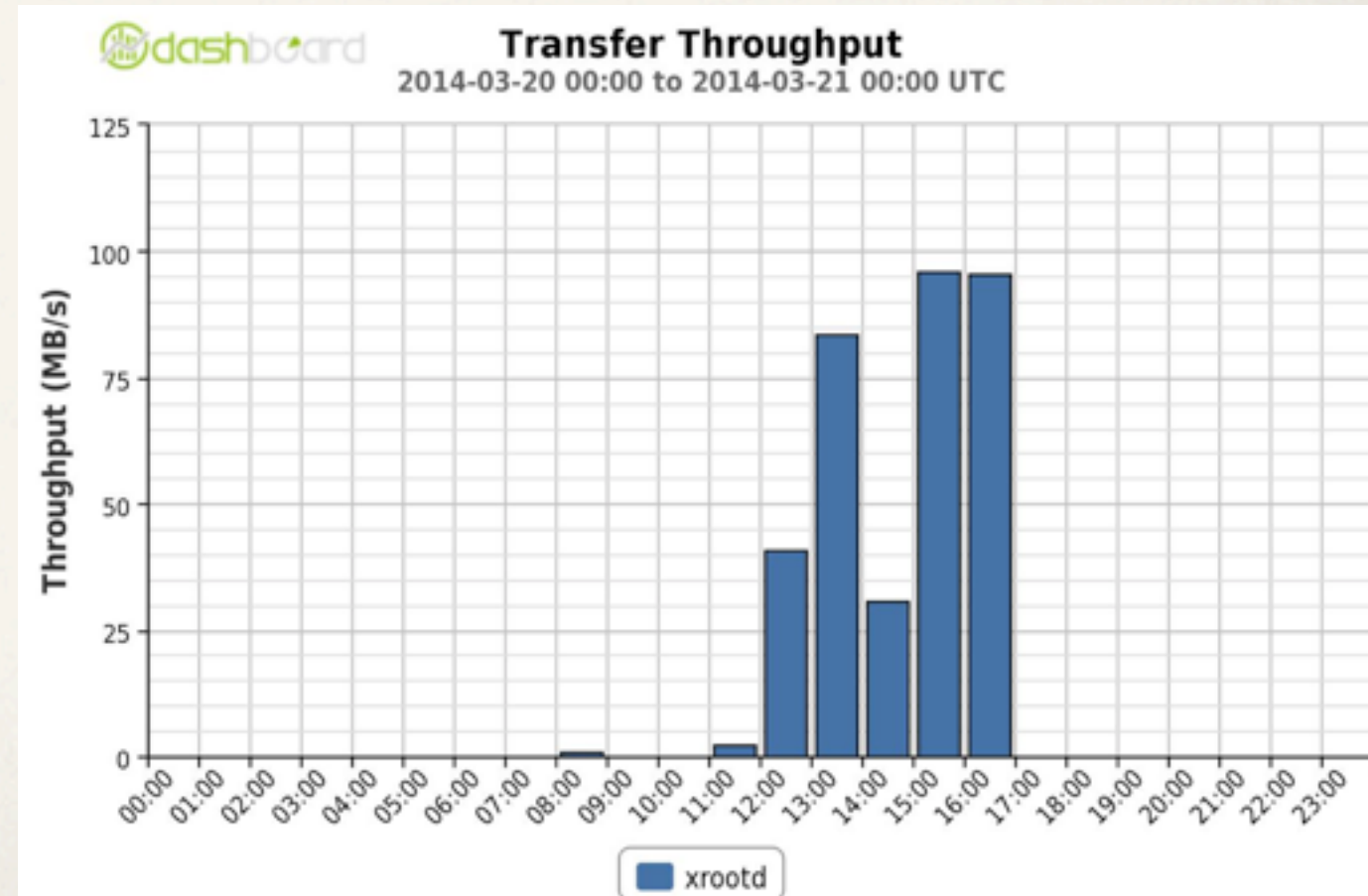1 Gig/s/Njobs
(and time increases)

NAT box can be upgraded but it's there for a reason ...

# Monitoring site activity



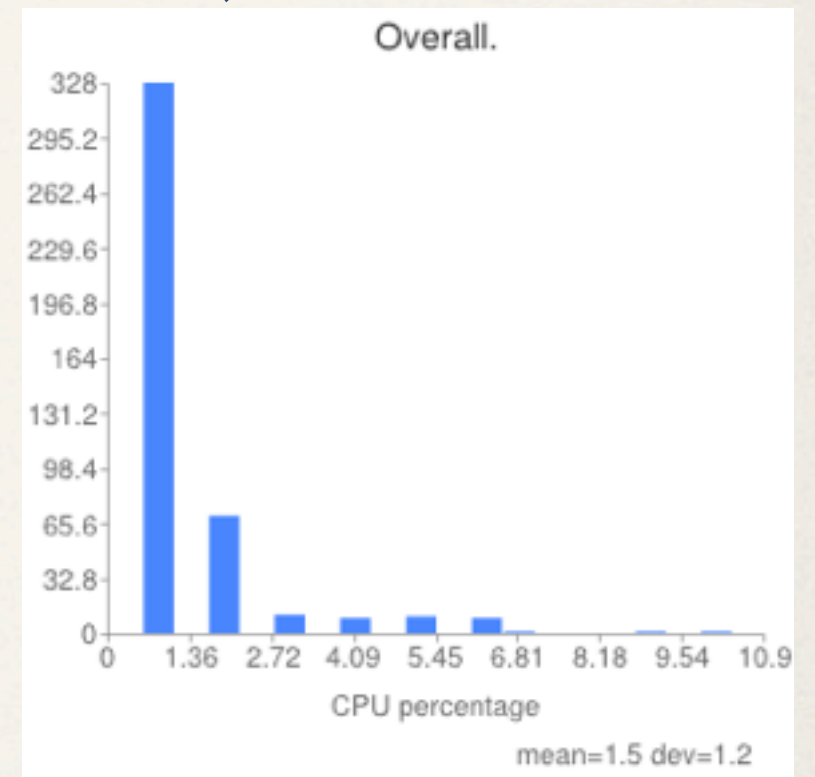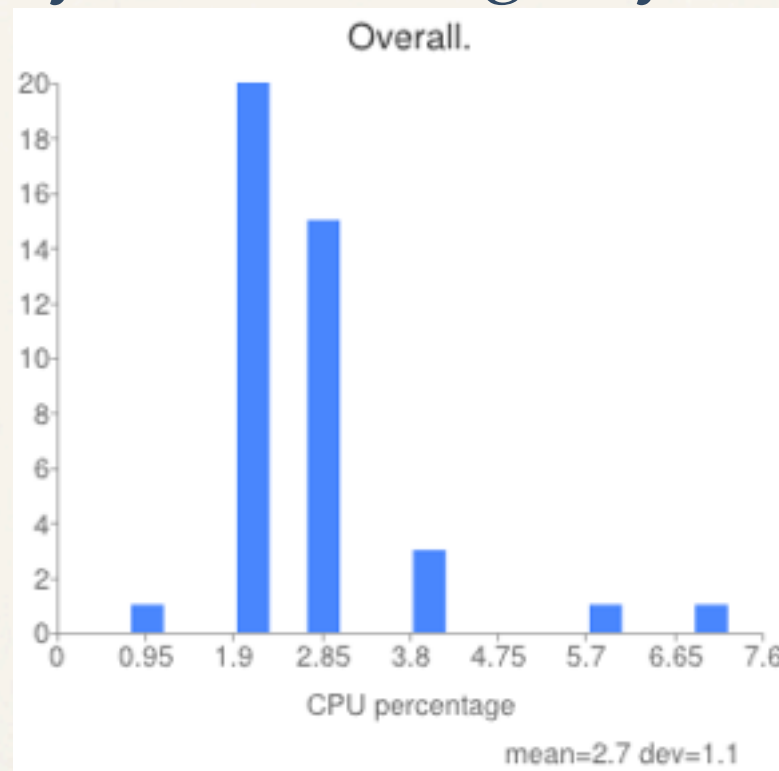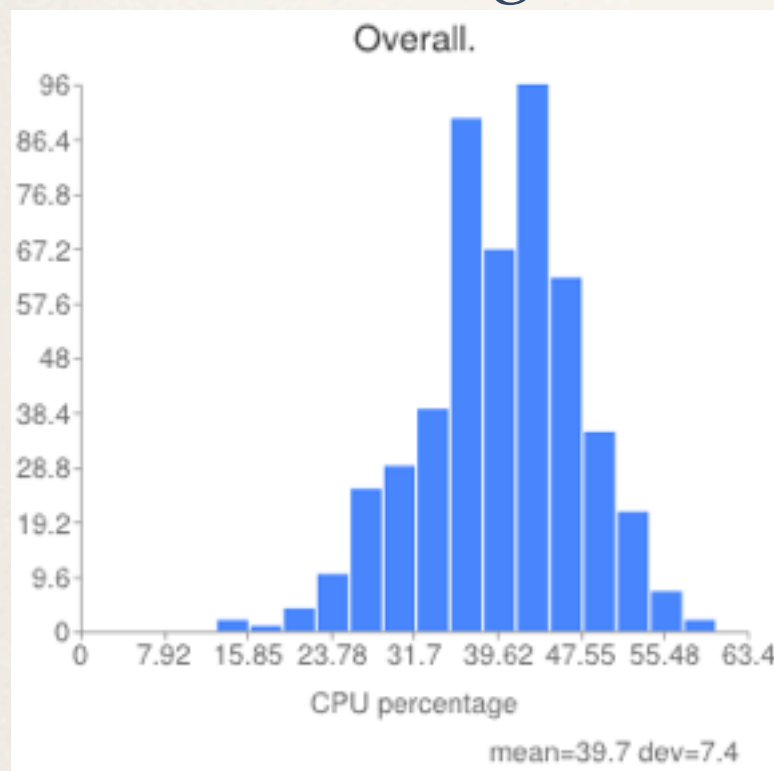Impact of those tests at ECDF switch
(Outward is blue)
Not an issue compared to FTS rate

* Easily seen in WLCG dashboard (isolating src: Lancs, dest:ECDF)

* This "bottleneck" is also a "bandwidth limit" so if removed we may need another: eg. Proxy server (not available for DPM sites by default) or xrootd Plugin (also not available by default)



**Transfer Throughput**
2014-03-20 00:00 to 2014-03-21 00:00 UTC

# Finding bottlenecks (2)...

✤ QMUL has a performant local Lustre setup with 10Gig to each WN from 100+ disk servers (see also backup slide).

✤ Xrootd though currently is via a single (yet untuned) server...



QM local panda test    QM - local but via xrootd    QM - from Oxford

✤ All this will be resolved, but shows it's worth testing...

# Atlas UK plans and initiatives

## Diskless Tier 2 at ANALY-UCL

* Limited local admin support  ideally wouldn't maintain storage..

* Decent WAN Links

* Initial H->WW tests work ⟶

* Soon ready for production : but how to broker jobs properly?

## Plans in Edinburgh:

* Opportunistic compute (industry clusters)

* Opportunistic storage ("RDF" (20PB non-HEP store))



Overall.

CPU percentage

mean=23.0 dev=7.3

# CMS - UK traffic

* Tier2 server sites for CMS are IC and RALPP (dCache), Brunel (DPM), others are clients.

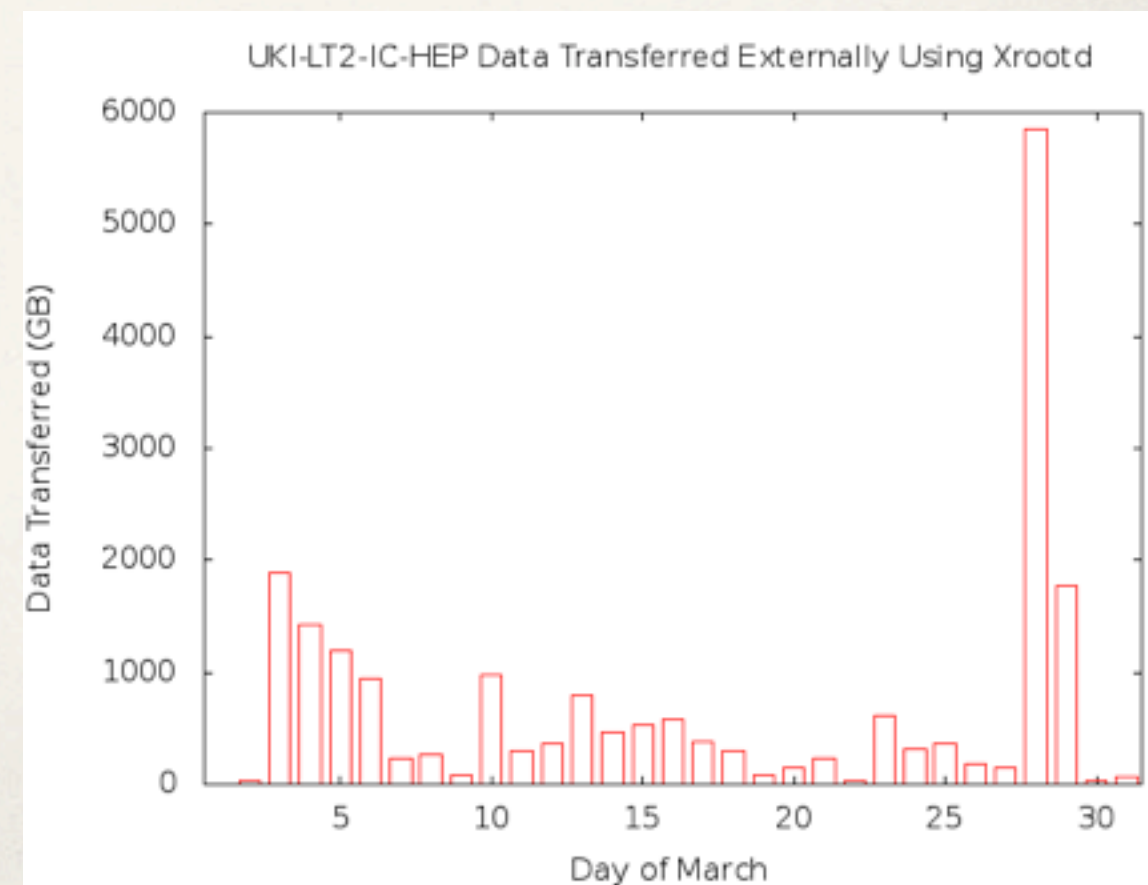* IC FTS/Xrootd traffic monitoring from billing logs

  * (Not in wlcg monitoring as they would prefer xrd monitoring plugin to be in dCache itself if it's required.)

  * gridftp out for March: 79 TiB

  * xrootd out for March: 21 TiB

  * Total WAN link traffic: 226 TiB (discrepancy users staging to home SE ?)



UKI-LT2-IC-HEP Data Transferred Externally Using GridFTP



UKI-LT2-IC-HEP Data Transferred Externally Using Xrootd

# CMS UK tests and plans

* Tests of analysis jobs performed ~2 years ago

* 5 different tests reading same file from client at Imperial - leftmost just reading (most I/O), rightmost is most "realistic"



Rebecca Lane - Imperial

---

"Diskless" plans for "DICE" Hadoop cluster in Bristol :

* Decent 10 Gbit/s (upgradable to 20 Gbit) link

* User jobs running on data at RALPP now - rigorous testing planned..

# A word on ATLAS / CMS comparisons

* From email chain and subsequent discussion with Brian in ROOT I/O w/g ...:

  * CMS tell sites to plan for 1MB/s per analysis job; usage tends to average 500KB/s.

  * CPU efficiency is around 75-80%.

* Not the same as ATLAS: e.g. "H->WW" code needs 20 MB/s to be 100%

* CMS "analysis" can involve reconstruction - higher CPU:

  * Its not a flaw or surprise if ATLAS get 10% eff where CMS have 100%

  * Not the same impact and issues in using federations.

* But of course ATLAS can and should improve their I/O, and have an opportunity with xAOD (Run2 format): the optimisation work has begun..

# Finally a word on HTTP

- ✤ Interest in HTTP from sites (mainly because its not HEP specific)

- ✤ Also from "small" experiments (again it's something they recognise..)

- ✤ Current uses (e.g. Rucio) are management ones (e.g. replace SRM...)

- ✤ But performance for data access also seems OK ..

(Single jobs, between ECDF (epel-test) and CERN (trunk) TEST boxes and using TWebFile not Davix )

# Conclusions

* Production level federations in the UK for ATLAS, CMS (and ALICE).

  * ATLAS and CMS cannot be directly compared and ATLAS will be limited (a bit) by current bandwidth (not necessarily a problem)

  * Opportunistic and diskless sites starting to be used.

* Starting to understanding bottlenecks

  * Monitoring important (good to see a discussion later...)

  * But also need to exercise (some) control ... (plugins or proxies..)

* Http/DAV will be used for ATLAS Rucio and "small"VOs: need to evaluate if its a reasonable alternative for the above use cases.
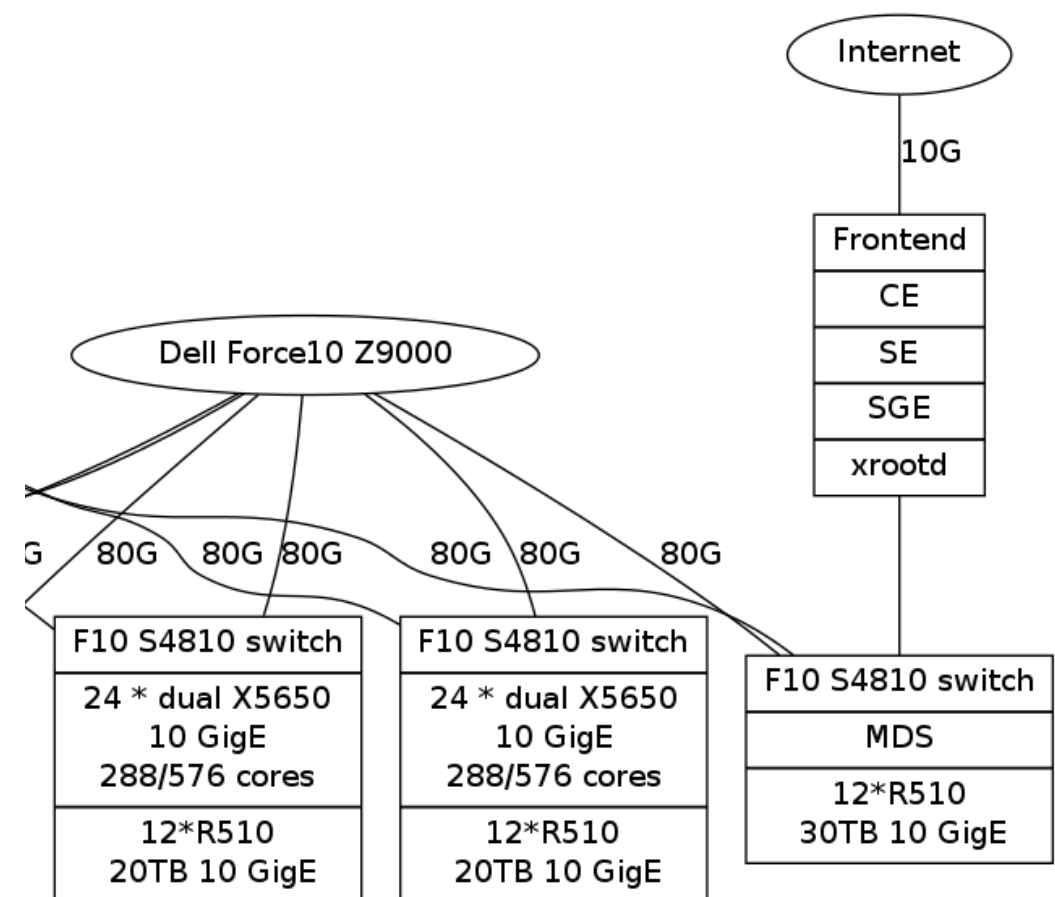
# Backups and background

# DPM

✤ UK been using xrootd/federations with DPM for a while

✤ DPM sites all have xrood redirector on the "headnode"  along with other services (e.g. SRM)

✤ Data transfer requests (local or remote) are redirected to disk server itself so transfers benefit from full bandwidth

✤ Almost all UK DPM sites now using xrootd for local atlas and cms traffic

# QMUL – Lustre and Storm

Highly optimized local access via Lustre will outperform xrootd access in current setup

- Local jobs reading from Lustre filesystem get the benefit of 10gig connections for WNs and the 100 disk servers

- Single xrootd server can only get 1/100 of this bandwidth to the servers.

- Xrootd not integrated with StoRM – r/o access for xrootd in atlas group.

- Traffic to WAN goes through a NAT.

Internet

10G

Frontend
CE
SE
SGE
xrootd

Dell Force10 Z9000

G   80G   80G  80G        80G   80G        80G

| F10 S4810 switch | F10 S4810 switch |
|---|---|
| 24 * dual X5650 10 GigE 288/576 cores | 24 * dual X5650 10 GigE 288/576 cores |
| 12*R510 20TB 10 GigE | 12*R510 20TB 10 GigE |

F10 S4810 switch
MDS
12*R510
30TB 10 GigE

QMUL Local: Chris Walker